

SOFTWARE

Open Access



ISMapper: identifying transposase insertion sites in bacterial genomes from short read sequence data

Jane Hawkey^{1,2*}, Mohammad Hamidian³, Ryan R. Wick¹, David J. Edwards¹, Helen Billman-Jacobe², Ruth M. Hall³ and Kathryn E. Holt¹

Abstract

Background: Insertion sequences (IS) are small transposable elements, commonly found in bacterial genomes. Identifying the location of IS in bacterial genomes can be useful for a variety of purposes including epidemiological tracking and predicting antibiotic resistance. However IS are commonly present in multiple copies in a single genome, which complicates genome assembly and the identification of IS insertion sites. Here we present ISMapper, a mapping-based tool for identification of the site and orientation of IS insertions in bacterial genomes, directly from paired-end short read data.

Results: ISMapper was validated using three types of short read data: (i) simulated reads from a variety of species, (ii) Illumina reads from 5 isolates for which finished genome sequences were available for comparison, and (iii) Illumina reads from 7 *Acinetobacter baumannii* isolates for which predicted IS locations were tested using PCR. A total of 20 genomes, including 13 species and 32 distinct IS, were used for validation. ISMapper correctly identified 97 % of known IS insertions in the analysis of simulated reads, and 98 % in real Illumina reads. Subsampling of real Illumina reads to lower depths indicated ISMapper was able to correctly detect insertions for average genome-wide read depths >20x, although read depths >50x were required to obtain confident calls that were highly-supported by evidence from reads. All ISAb1 insertions identified by ISMapper in the *A. baumannii* genomes were confirmed by PCR. In each *A. baumannii* genome, ISMapper successfully identified an IS insertion upstream of the *ampC* beta-lactamase that could explain phenotypic resistance to third-generation cephalosporins. The utility of ISMapper was further demonstrated by profiling genome-wide IS6110 insertions in 138 publicly available *Mycobacterium tuberculosis* genomes, revealing lineage-specific insertions and multiple insertion hotspots.

Conclusions: ISMapper provides a rapid and robust method for identifying IS insertion sites directly from short read data, with a high degree of accuracy demonstrated across a wide range of bacteria.

Keywords: Insertion sequence (IS), Bacteria, Genomics, Short read analysis, Tuberculosis, Antimicrobial resistance

Background

An insertion sequence (IS) is a small transposable element that encodes the proteins required for its own transposition. The ISfinder database [1] currently contains over 500 distinct IS. During transposition some ISs create

direct repeats, or target site duplications, in the sequences into which they are integrating. The presence and length of these duplications vary widely between ISs and are characteristic of individual IS [2]. Rates of transposition vary between ISs and host species, but are frequently in the order of the rate of nucleotide substitutions, making IS activity one of the more dynamic evolutionary forces at play in many bacterial genomes. The movement of ISs can also have functional consequences for bacterial genomes. ISs have been implicated in large changes to genome structure, by expanding in copy number in microbial

* Correspondence: hawkey.jane@gmail.com

¹Department of Biochemistry and Molecular Biology, Bio21 Molecular Science and Biotechnology Institute, The University of Melbourne, Parkville, VIC 3010, Australia

²Faculty of Veterinary and Agricultural Science, The University of Melbourne, Parkville, VIC 3010, Australia

Full list of author information is available at the end of the article

genomes, with subsequent loss of ISs resulting in inactivation of genes, pseudogene formation, mediating deletion of intervening sequences between two copies of the IS, or rearrangements of the genome [3].

In addition, IS insertions upstream of protein coding sequences can result in their enhanced expression, leading to different phenotypes depending on the function of the over-expressed gene. There are several known examples of IS-mediated gene expression leading to clinically important increases in antimicrobial resistance. For example, increased resistance to fluoroquinolones such as ciprofloxacin can result from the insertion of IS1 or IS10 upstream of the *acrEF* efflux pump in *Salmonella* Typhimurium [4], or the insertion of IS186 upstream of the *acrAB* efflux pump in *Escherichia coli* [5]. In *Acinetobacter baumannii*, insertion of ISAbal or ISAbal25 upstream of the intrinsic beta-lactamase *ampC* can cause resistance to third generation cephalosporins including ceftazidime and cefotaxime [6, 7]. Insertions of the same IS in nearby locations can generate a composite transposon, capable of mobilizing the intervening sequence and transferring it to new genomic locations. For example, the composite transposon Tn6168 was generated spontaneously via insertions of ISAbal on either side of *ampC*, including one copy of ISAbal that upregulates *ampC* expression [8]. Tn6168 has then transferred into different *A. baumannii* backgrounds, conferring horizontally-acquired resistance to third generation cephalosporins [8].

IS insertions also result in the upregulation of virulence genes in clinically important human pathogens. For example, an outbreak of tuberculosis in Spain in the 1990s was associated with the B strain of *Mycobacterium bovis* carrying an insertion of IS6110 in the promoter region of the virulence gene *phoP*, resulting in its upregulation [9]. In *Neisseria meningitidis*, insertion of IS1301 in the middle of the capsule locus has been shown to cause increased expression of operons on either side of the IS, contributing to protection from the human immune system and enhanced pathogenicity [10]. ISs have also been shown to enhance niche adaptation in bacteria, for example IS1247 insertion upstream of *dhlB* in *Xanthobacter autotrophicus* results in increased resistance to bromoacetate [11]. This region has also been mobilised by the IS and transferred to a plasmid [11]. In *E. coli*, IS3 has been shown to up-regulate threonine expression, allowing the bacteria to adapt to a low-carbon environment and utilise threonine as its sole carbon source [12].

The profiling of IS insertion patterns has been used for typing purposes in numerous bacterial species of importance to human health. For example, copy number and position of IS200 in *Salmonella enterica* [13], IS6110 in *Mycobacterium tuberculosis* [14], IS1004 in

Vibrio cholerae [15] and ISAbal in *A. baumannii* [16] has been used to profile these bacterial pathogens, allowing the identification and tracking of distinct subtypes. To date, IS-based typing schemes for various bacteria have relied on digesting the genome followed by either hybridizing IS probes to fragments in a gel or PCR probing [13–15]. The detection of precise insertion sites can be achieved using PCR, and may be done for typing purposes [17] or for the detection of functionally important insertions [7, 9].

With the advent of cheap high-throughput short-read sequencing, whole genome sequencing (WGS) of bacteria is increasingly common and is replacing traditional methods for characterizing and typing bacterial genomes. Unfortunately the detection of ISs is complicated wherever read lengths are shorter than the length of the IS, as is the case for platforms that are currently most widely used – Illumina and Ion Torrent. IS insertion sites can readily be identified in finished bacterial genomes or in draft assemblies of genomes with single-copy ISs, using tools such as nucleotide BLAST or ISfinder [1]. However where multiple copies of the same IS are present within a single genome (including on the chromosome and/or plasmids), this complicates assembly of short-read data and makes IS insertion sites difficult to identify reliably. The IS detection problem can be resolved using long-read sequencing technologies such as the SMRT Cell (Pacific Biosciences) or MinION (Oxford Nanopore) platforms; however given the relative cost efficiency and reliability of short-read sequencing, together with the current widespread use of Illumina for bacterial WGS and wealth of available short-read data for clinically important bacteria, there remains a need for a simple tool to identify IS insertion sites from short-read data.

Several studies report the use of mapping-based approaches to identify IS insertion sites from bacterial short-read data [18, 19], however none provide software code or validation of the approach used. There are tools available for detecting transposons or structural variation in genomes, for example MindTheGap [20], BreakDancer [21], and Mobster [22], however these do not perform well in the identification of IS in bacterial genomes nor were they designed to do so. Some programs could potentially be used for this purpose, such as Relocate [23] and RetroSeq [24], however these require additional input or prior knowledge about the IS which may not always be available. TIF (Transposon Insertion Finder) [25] and *breseq* [26] could potentially be used for the detection of IS insertion sites in bacterial genomes, however they were not designed specifically for this purpose and did not perform well on our data sets (see Results).

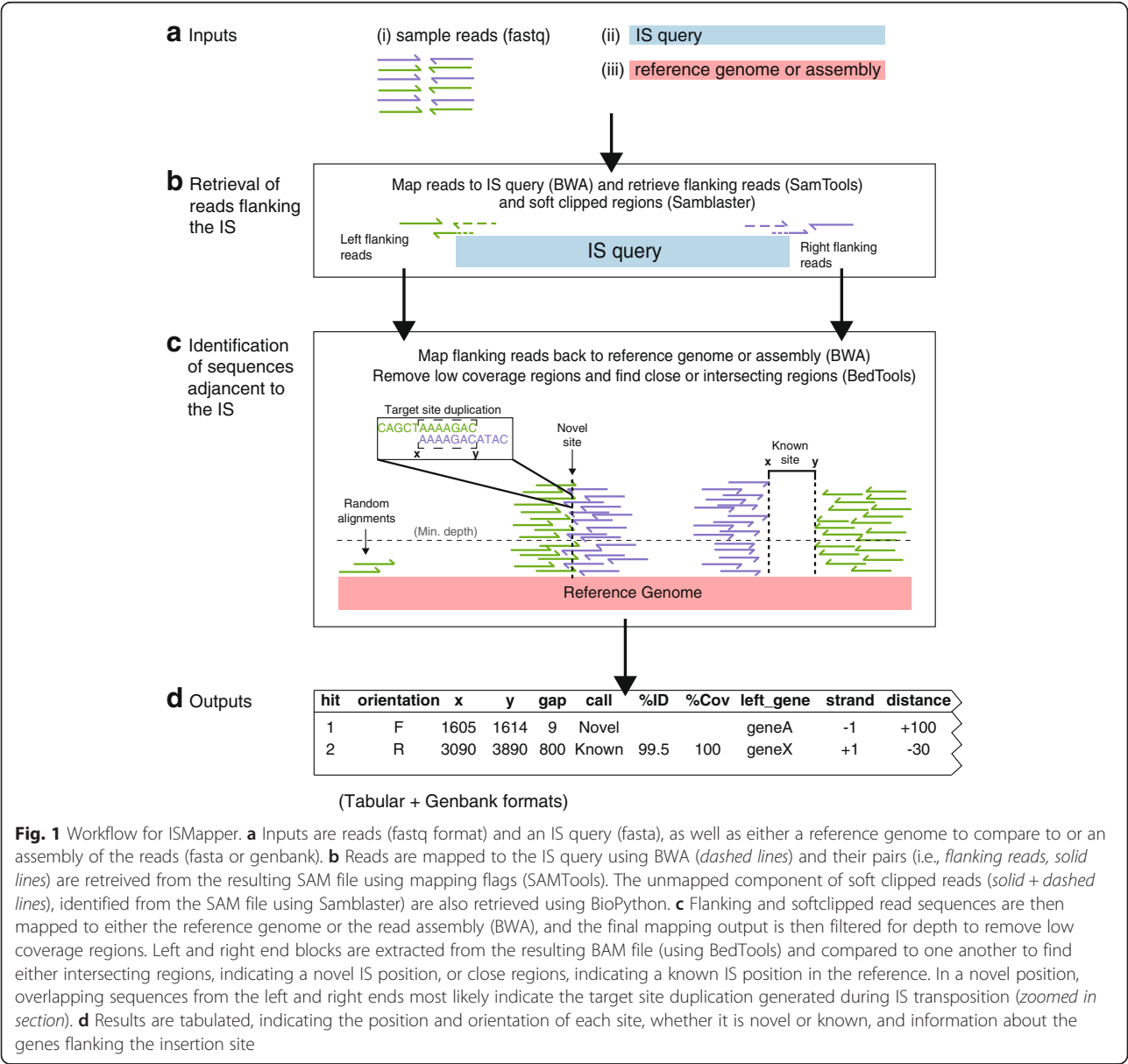
Here we present a rapid and robust tool for accurate detection of ISs, including insertion site and orientation,

direct from short-read data. The method is freely available in the form of open-source code called ISMapper, and here we validate its use via analysis of simulated and real short-read data from a range of ISs and bacterial species. ISMapper requires short reads and query IS sequences as input, and can be used either for typing against a reference genome or to assist with manual resolution of complex short-read assemblies.

Implementation

An overview of the ISMapper workflow is shown in Fig. 1. ISMapper takes as input: (i) a set of paired end Illumina reads for an isolate of interest, (ii) an IS query sequence in fasta format, and (iii) either a reference genome (for typing) or an assembly of the read set (for

assembly improvement), in GenBank or FASTA format (Fig. 1a). Paired end Illumina reads are mapped to the IS query sequence using BWA-MEM (v0.7.5a or later) [27]. From the resulting alignment file (SAM format), unmapped reads whose pairs map to the end of the IS query sequence (that is, reads representing the sequences directly flanking the IS) are extracted using SAMtools view (v0.1.19 or later) [28] to retrieve reads based on SAM flags (Fig. 1b). Specifically, left flanking reads (taking input sequence as left to right) are extracted using flag ‘-f 36’ and right flanking reads are extracted using flag ‘-F 40 -f 4’ and stored in separate BAM files, which are then converted to FASTQ format using BedTools (v2.20.1) [29]. In addition, Samblaster (v0.1.21) [30] is used to extract from the SAM file any

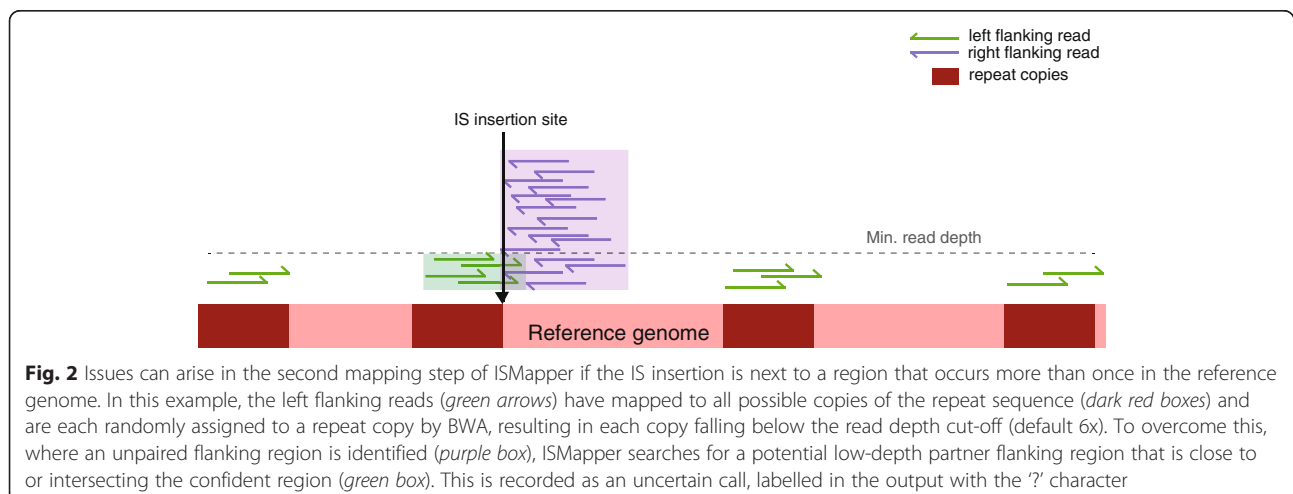


reads that map to the end of the IS and extend into the neighbouring sequence (i.e., “soft clipped” reads, Fig. 1b). The resulting FASTQ file is filtered using BioPython to extract the soft clipped portion of reads, where those sequences fit a specified size range (default 5–30 bp). The resulting sequences are sorted into left and right flanking sequences; these are each mapped separately to the reference genome or assembly using BWA-MEM, to identify the location(s) of the query IS in the genome under analysis (Fig. 1c). Insertion site information is extracted from the resulting alignments using BedTools (coverage command) to summarise coverage of the reference by left and right flanking reads; these are filtered by read depth (default, minimum read depth $\geq 6x$) to minimize false positive hits, and regions that overlap or are separated by a short distance (default, ≤ 100 bp) are merged using BedTools (merge command). Pairs of left and right flanking regions that likely represent either side of the same IS insertion are identified on the basis of positional information, using BedTools (intersect and closest commands). Left and right regions that overlap are considered to indicate a novel IS insertion not present in the reference, with the overlap resulting from target site duplication arising during IS transposition (Fig. 1c, novel site). Coordinate x refers to the left side of the target site duplication, and y refers to the coordinate of the right side of the target site duplication. In cases where the left and right flanking regions are extremely close but do not overlap, x refers to the inner end of the left-most region, and y refers to the inner end of the right-most region (as per a known site, see below). Where left and right regions are separated by a sequence that is approximately the length of the IS query, the intervening sequence is extracted and compared to the IS query using nucleotide BLAST+ (v2.2.25 or later) [31] to confirm whether this is a known insertion site that is present in the reference (Fig. 1c, known site). In this case, coordinate x refers to

the inner end of the left-most block, and y refers to the inner end of the right-most block. Coordinate x is always the smallest number and y always the largest, regardless of IS orientation.

Extensive testing of ISMapper revealed that it was sometimes unable to resolve IS positions that were adjacent to a repeat region (segments of DNA that were repeated multiple times around the genome; see Results). This is because when the IS-flanking reads were mapped back to the reference genome, those that belonged to the neighbouring multi-copy sequence were randomly assigned by BWA-MEM to the various locations of the repeat sequence, resulting in low read depth at the ‘true’ IS-adjacent copy of the multi-copy sequence, which can fall below the minimum depth filter (Fig. 2). In such cases, the sequence on the other side of the IS is usually not a multi-copy sequence and thus does not suffer the same problem, and so is usually identified as a confident IS-flanking region without a corresponding partner region (Fig. 2, purple block). Therefore, when ISMapper identifies an un-partnered IS-flanking region, it checks the original alignments for evidence of a nearby low-coverage partner region that failed to pass the depth filter and returns this as a potential but uncertain IS location, indicated by a ‘?’ character in the results table (see example of low-coverage partner region in Fig. 2, green block).

ISMapper generates two main output files summarizing the results: (i) a GenBank file of the reference sequence, annotated with the IS-flanking regions and (ii) a table indicating the locations and characteristics of each IS-flanking region identified (Fig. 1d). The table includes details of the location of the IS insertions (indicated by coordinates x and y); the distance between the left and right flanking regions (where a negative number indicates an overlap of left and right regions, indicating the size and sequence of the target site duplication); a call as



to whether the insertion is present in the reference or is a novel insertion site (and, where the insertion site is present in the reference, the percent coverage and sequence homology with the IS query); and details of the gene(s) closest to the IS insertion site (including locus tag, product, gene name and distance from the IS to the start codon). Insertions are also marked to indicate less confident calls. A '*' indicates an imprecise hit; i.e., where the gap between left and right regions is larger than expected for a novel insertion, but is not consistent with an IS insertion at that location in the reference. A '?' indicates an uncertain hit, where only one end (left or right of the predicted insertion) passes the minimum read depth threshold; this often occurs when the IS is inserted within or adjacent to a multi-copy sequence, as described above (Fig. 2). When run in assembly improvement mode, the table produced is simpler and indicates which contigs are predicted to end adjacent to the IS (indicating left or right orientation), assisting the user to decide whether some contigs could be joined together based on the available IS evidence.

ISMMapper is lightweight code – a test run on a laptop computer (MacBook Air) with 8GB of RAM and a 1.3GHz i5 processor was able to analyse a read set comprising 2.5 million 100 bp paired-end reads in approximately ten minutes for a single IS query. Because ISMapper analyses each read set and query IS independently, screening of multiple read sets and query IS can be easily performed in parallel across multiple cores. To facilitate easy compilation of results from multiple jobs, ISMapper includes a Python script to cross-tabulate results from multiple read sets, generating a single summary table per query IS (script 'compiled_table.py').

Results and discussion

Validation of IS detection using simulated reads

Nine publicly available finished genomes from a variety of bacterial genera, and including both chromosomes and plasmids, were downloaded from NCBI (Table 1). ISfinder [1] was used to identify the ISs present in each finished genome sequence. All sequences that had >50 % identity to a sequence in ISfinder and were present in at least two copies were tested using ISMapper (with the query IS being sourced from curated references in ISfinder). Nucleotide BLAST+ was used to confirm the precise locations and orientations for each query IS in all genomes (total 251 insertions of 17 distinct IS, see Table 1). Short reads (100 bp) were simulated from each genome sequence using the wgsim command in SAMtools (v0.1.19), with default parameter settings.

ISMMapper was run with default parameter settings on each combination of genome, query IS and simulated reads. ISMapper was able to accurately locate each IS position and its orientation (ranging between 2 and 61

Table 1 Validation of ISMapper using simulated reads

Isolate	Accession	IS	Found	Orientation
S. Typhi CT18	NC_003198	IS200	26/26	26/26
		IS1	3/3	3/3
S. Typhimurium LT2	NC_003197	IS200	6/6	6/6
		ISSty2 ^{b (1)}	2/2	2/2
		IS1351 ^{b (2)}	2/2	2/2
S. Typhi plasmid pHCM1	NC_003384	IS26	4/4	4/4
		ISVsa5	2/2	2/2
		IS1	5/5	5/5
S. Paratyphi plasmid pAKU_1	AM412236	IS26	5/5	5/5
		ISVsa5	2/2	2/2
		IS1	7/7	7/7
<i>K. pneumoniae</i> plasmid pNDMAR	JN420336	IS3000	3/3	3/3
		IS26 ^a	3/5	3/5
		ISEcp1	2/2	2/2
<i>Yersinia pestis</i> C092	NC_003143	IS100 ^a	43/44	43/44
		IS1661 ^{a, b (1)}	7/8	7/8
		IS1541 ^a	61/64	61/64
<i>Escherichia coli</i> O104:H4	NC_018658	IS1	10/10	10/10
		IS421	4/4	4/4
		IS609	4/4	4/4
		ISEc1	4/4	4/4
		ISKpn26	4/4	4/4
<i>E. coli</i> O157:H7	NC_002695	IS629 ^{a, b (1)}	17/18	17/18
		IS609	2/2	2/2
		ISEc1 ^{b (1)}	4/4	4/4
		ISEc8 ^a	9/10	9/10

^aIndicates ISMapper was unable to resolve some IS positions due to repeat regions

^bIndicates ISMapper incorrectly identified an insertion site, the number of these sites are indicated in brackets

positions per genome) for the majority of genomes (Table 1). In total, 97 % of IS insertions were correctly detected, with a false positive rate of 2.1 % ($n = 6$). The exceptions occurred in three genomes (*K. pneumoniae* plasmid pNDMAR, *Y. pestis* C092 and *E. coli* O157:H7), in which ISMapper correctly identified 151 IS insertion sites and failed to identify nine (94 % detection). Closer inspection revealed that the missed IS were each located next to multi-copy repeat sequences, complicating the second mapping step as discussed above and outlined in Fig. 2. Switching on reporting of all alignments above a mapping score threshold of 30 (–a and –T 30 in BWA-MEM) enabled the detection of a further IS100 site in *Y. pestis*. By default this option is turned off in ISMapper as it tends to create noise in the mapping, making it more difficult to distinguish true and false positives;

however this can be useful if an IS site of interest is known or suspected to be flanked by further repeats.

Validation of IS detection using real Illumina read sets derived from isolates with finished genomes

Next we validated ISMapper using six finished genomes for which both Illumina read data and finished genomes were publicly available (Table 2). Each finished genome sequence was analysed with ISfinder [1] to identify query ISs for testing as described above, and nucleotide BLAST was used to confirm the precise locations and orientations of each IS in each genome. The resulting test set comprised 106 insertions of 14 query ISs. Using default settings, ISMapper was able to accurately identify each IS insertion site and its orientation, between 2 and 26 per genome, for the majority of genomes (Table 2). In total, 104 (98 %) IS insertions were correctly detected by ISMapper, with 3 IS insertions falsely detected (false positive rate of 2.5 %, $n = 3$). Three of four IS431mec insertions in *Staphylococcus aureus* TW20 were correctly detected, however the fourth was missed by ISMapper as it was flanked by another IS431mec and further repeat sequences. Two of three IS1 insertions in *Salmonella* Typhi CT18 were correctly detected however a third, located between *tviE* and *tviD*, was problematic. ISMapper identified the region flanking the IS at *tviE*, but did not detect any corresponding region in *tviD*. To investigate why, we mapped the entire Illumina read set to the CT18 chromosome reference sequence, which showed that there were no reads derived from the region between *tviD* to *tviA*. Therefore this region appears to have been deleted during culture in the laboratory prior

to the extraction of DNA for Illumina sequencing. This region encodes the biosynthesis of the Vi capsule of *S. Typhi*, and is known to be lost sporadically during culture [32]. This illustrates that situations where one end of the IS is detected but the other is not can often be 'accurate' in the sense that the result reflects underlying structural variation in the genome, including potentially IS-mediated deletions.

Detection of antibiotic resistance-mediating IS insertions in *Acinetobacter baumannii*, confirmed by PCR

The genomes of seven ceftazidime resistant *A. baumannii* isolates, belonging to global clone (GC) 1, were sequenced via Illumina HiSeq to generate 100 bp paired end reads. Resistance gene screening of the Illumina data using SRST2 [33] and the ARG-Annot database [34] confirmed earlier PCR data indicating that none of these isolates carried acquired extended spectrum beta-lactamase (ESBL) genes that can confer resistance to third-generation cephalosporins. However, it is known that the insertion of ISAbal1 upstream of the intrinsic chromosomally encoded *ampC* beta-lactamase gene can cause increased resistance to third-generation cephalosporins in *A. baumannii* [6].

We used ISMapper to screen for the ISAbal1 query sequence (accession AY758396), sourced from ISfinder [1]. Using default parameters, ISMapper identified ISAbal1 insertions in all seven GC1 genomes. IS positions were assessed relative to the finished genome sequence of *A. baumannii* GC1 reference A1 (accession CP010781). ISMapper found between 3 and 5 ISAbal1 insertions in each GC1 isolate, including an insertion upstream of

Table 2 Validation of ISMapper using real Illumina reads for which finished genomes were also available

Isolate	Genome accession	FastQ accession	IS	Found	Orientation
<i>Streptococcus suis</i> P1/7	AM946016	ERR225612	ISSu3	4/4	4/4
			ISSu4 ^{b (2)}	2/2	2/2
<i>Staphylococcus aureus</i> TW20	NC_017331	ERR043367	ISep3	3/3	3/3
			IS256	8/8	8/8
			IS431mec ^a	3/4	3/4
			IS1181	2/2	2/2
<i>Klebsiella pneumoniae</i> NJST258_1	CP006923	SRR1166975	ISKpn1	5/5	5/5
			IS5B	8/8	8/8
			IS903B ^{b (1)}	2/2	2/2
			IS1294	3/3	3/3
			ISKpn18	2/2	2/2
			ISKpn26	7/7	7/7
<i>S. Typhi</i> CT18	NC_003198	ERR343331	IS200	26/26	26/26
			IS1 ^a	2/3	2/3
<i>S. Typhi</i> Ty2	AE014613	ERR343332	IS200	26/26	26/26

^aindicates ISMapper was unable to resolve some positions due to repeat regions

^bindicates ISMapper incorrectly identified an insertion site, the number of these sites are indicated in brackets

ampC in all 7 genomes that was in the orientation required to induce upregulation and explain the observed cephalosporin resistance phenotype (Fig. 3). In addition, out of 29 total ISAbal insertions, ISMapper was able to correctly identify 26 target site duplications (9 bp in the case of ISAbal). All ISAbal insertions were novel compared to the reference genome A1 (Fig. 3) and were confirmed using PCR, as described in [6].

Impact of read depth on ISMapper performance

To test the effect of read depth on the performance of ISMapper, each of the seven GC1 *A. baumannii* read sets were randomly subsampled to depths of approximately 10x, 15x, 20x, 25x, 50x, 75x and 100x, with ten replicates per depth level per read set. ISMapper was then run using default settings to screen for ISAbal insertions. The results indicated that at mean genome-wide read depths of approximately 20x, ISMapper was able to identify 95 % of insertions correctly (Fig. 4). However, all of these calls were either imprecise (gap size larger than expected) or uncertain (high coverage end paired with a low coverage end). An average genome-wide read depth of ~50x was required to find

all insertions, with confident calls for >60 %, however there was clearly some variation depending on read quality (Fig. 4). To achieve 100 % detection with high confidence, average genome-wide read depths of >75x were required (Fig. 4).

Comparison of ISMapper with TIF and breseq

The seven *A. baumannii* GC1 genomes were used to test both *breseq* [26] and TIF (Transposon Insertion Finder) [25]. *breseq* uses split read mapping to a reference genome along with statistical models to determine new junctions and deletions in the isolates of interest. As input, *breseq* takes paired end reads in FASTQ format, and a reference genome in Genbank format. The *breseq* manual indicates that new insertions of mobile elements can be determined by looking for 'JC JC' evidence types in the final html output. All seven *A. baumannii* isolates were screened using default parameters and the reference genome A1 (accession CP010781). In all cases, *breseq* was unable to identify any mobile element insertions, including no structural variation at the known ISAbal insertion sites, although many other types of structural variation were detected.

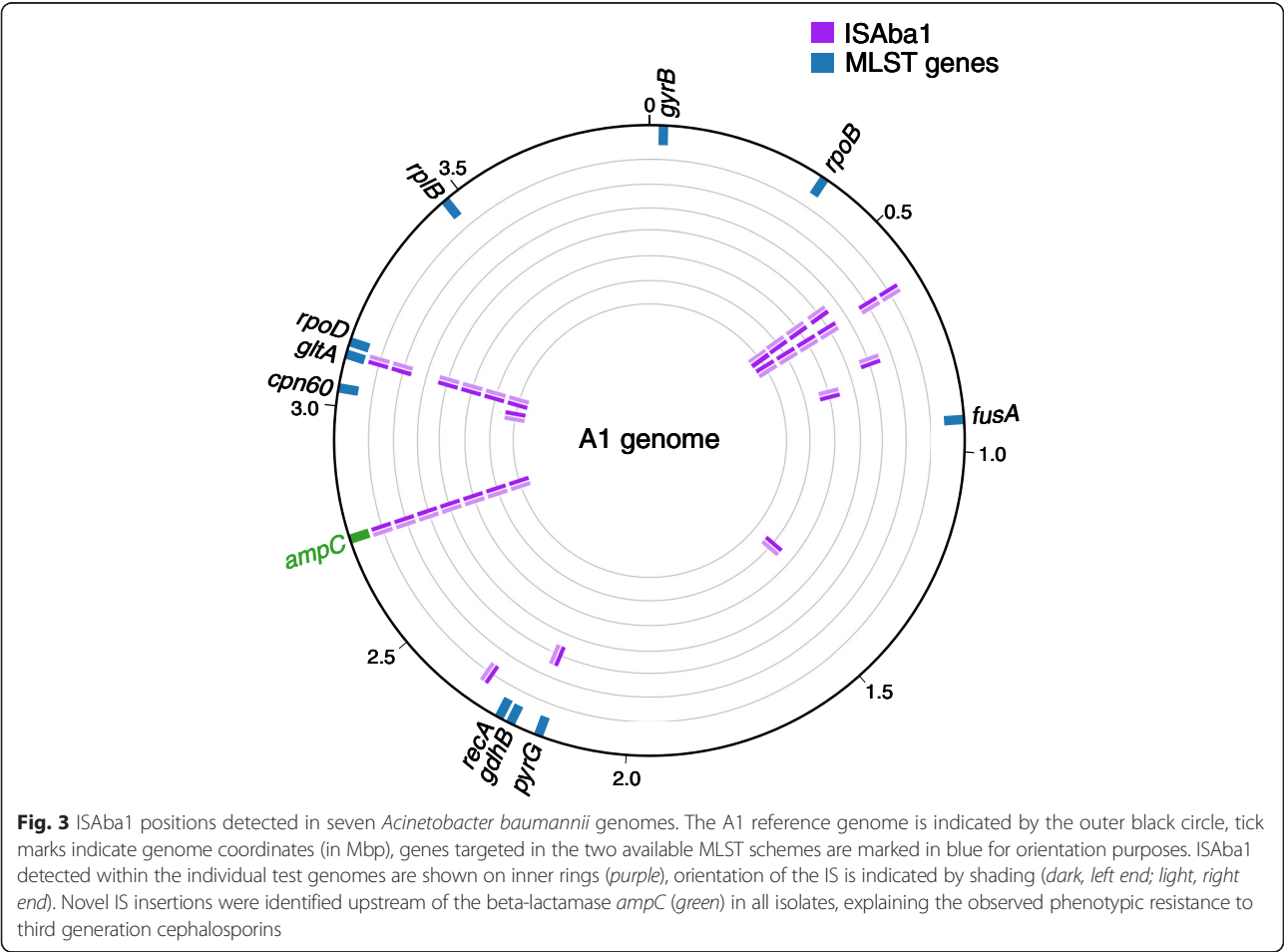


Fig. 3 ISAbal positions detected in seven *Acinetobacter baumannii* genomes. The A1 reference genome is indicated by the outer black circle, tick marks indicate genome coordinates (in Mbp), genes targeted in the two available MLST schemes are marked in blue for orientation purposes. ISAbal detected within the individual test genomes are shown on inner rings (purple), orientation of the IS is indicated by shading (dark, left end; light, right end). Novel IS insertions were identified upstream of the beta-lactamase *ampC* (green) in all isolates, explaining the observed phenotypic resistance to third generation cephalosporins

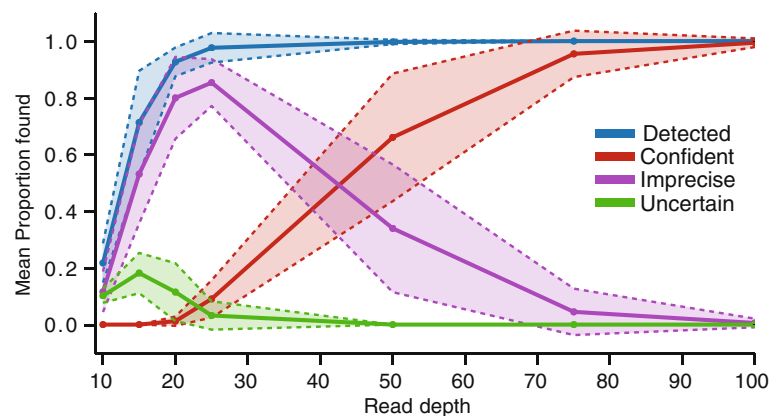


Fig. 4 ISAbal detection rate as a function of read depth, for seven *A. baumannii* GC1 genomes sequenced with Illumina. Mean (lines) and standard deviation (shaded areas) proportions of IS insertions correctly detected per genome, amongst 70 replicate read sets sampled at each depth level (7 read sets x 10 replicates each at read depths 10x, 15x, 20x, 25x, 50x, 75x, 100x). Blue, IS insertion site detected; red, detected with high confidence; purple, detected with low precision (larger than expected gap size between left and right flanking regions, indicated with '*' in output); green, detected with low confidence (high-depth evidence for one side only, low-depth evidence for the other, indicated with '?' in output)

TIF requires as input paired end reads (FASTQ format), the head and tail sequences (approximately 17 bp) of the IS of interest as well as the size of the target site duplication the IS makes during transposition. TIF uses regular expressions to search for the head and tail sequences in the reads, and these reads are then extracted and grouped by their target site duplications. Unfortunately, following communication with the authors, we were unable to get TIF to output any results using our data. Other disadvantages of TIF are the requirements to (i) specify the size of target site duplications (which not all IS make and is not always known), (ii) manually extract subsequences of the IS rather than inputting the complete sequence, and (iii) manually edit a Perl script in order to specify inputs to the program.

Example use case: exploration of IS6110 insertions in *Mycobacterium tuberculosis*

While IS insertions are thought to be important for shaping the evolution of bacteria in a variety of ways, high-resolution comparative genomic studies of bacterial pathogens have largely ignored ISs due to the difficulties associated with accurate detection of insertion sites from high-throughput short read data. An important example is IS6110 in *M. tuberculosis* [35]. Profiling of IS6110 insertions using PCR and restriction fragment based polymorphism (RFLP) based methods has been reported for typing purposes [36], and specific insertions have been linked to clinically relevant changes in function including in outbreak strains [9, 37, 38]. However while numerous studies have reported the genomic analysis of hundreds of *M. tuberculosis* isolates sequenced on the Illumina platform, these have not included analysis of IS6110 insertions. Thus, to demonstrate the utility of

ISMMapper for comparative profiling of ISs in an important bacterial pathogen, we analysed the distribution of IS6110 within 138 publicly available genomes representing the major lineages of *M. tuberculosis* [39]. Paired-end Illumina reads were downloaded from NCBI (ERP001731). A core genome phylogeny was generated from these reads by SNP (single nucleotide polymorphism) calling against reference genome H37Rv (accession NC_000962) (methods as described in [40]), followed by maximum likelihood phylogenetic inference on the SNP alignment using RAxML (GTR + G substitution model, 1000 bootstraps) to build a genome-wide phylogenetic tree. ISMapper was run with default settings to screen for insertions of IS6110 (accession X17348) in each read set, relative to reference genome H37Rv.

A total of 392 unique IS6110 insertion sites were identified by ISMapper, approximately one per 10 kbp of the 4.4 Mbp reference genome. The frequency of each insertion within each of the six main global lineages is shown in Fig. 5b. The data indicate multiple lineage-specific IS6110 insertions in lineages 2–6, but none that were shared by multiple lineages, suggesting that IS6110 insertions began to accumulate only after *M. tuberculosis* diverged into these distinct lineages. Isolates in the “modern” lineages 2–4 and in the West African lineage 5 had more IS6110 insertions overall, with far fewer insertions observed in the “ancient” East and West African lineages 1 and 6 (Fig. 5c). Lineage 2, which includes the highly successful Beijing sublineage, had the highest number of IS6110 although it was not the most common lineage in the collection ($n = 23$); it could be that these insertions contribute to the adaptive fitness of the Beijing lineage.

The spatial distribution of unique IS6110 insertions within the *M. tuberculosis* genome (Fig. 5d) revealed

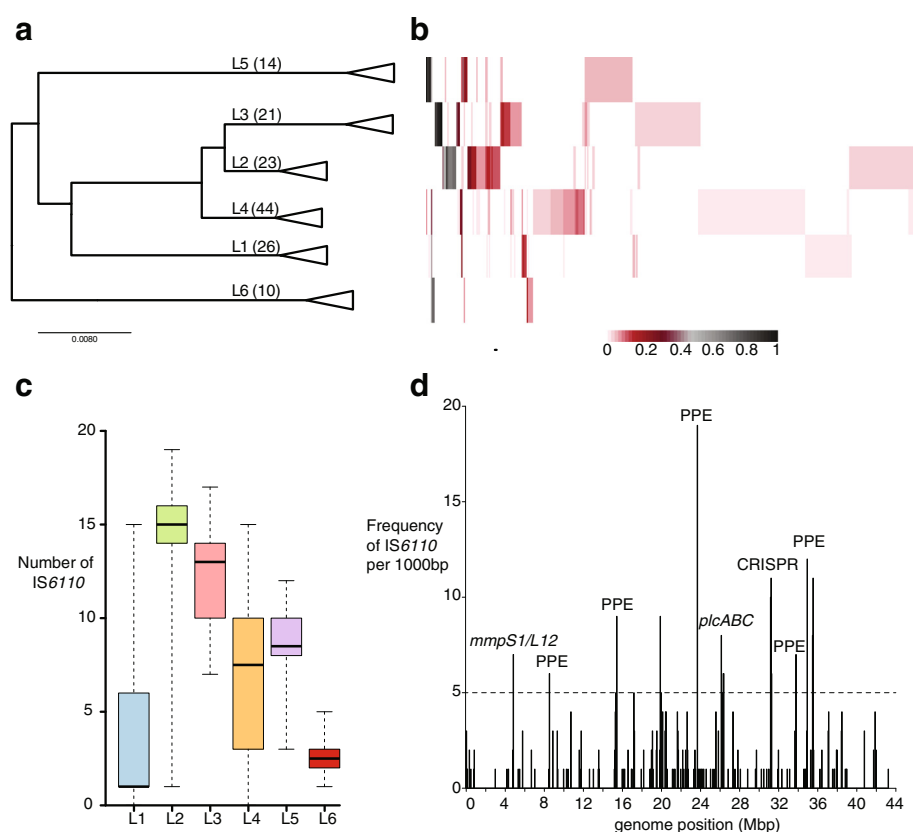


Fig. 5 Analysis of IS6110 insertions in a diverse set of *M. tuberculosis* genomes. Analyses are based on ISMapper analysis of publicly available Illumina paired end data from 138 genomes. **a** Phylogenetic tree for *M. tuberculosis* based on genome-wide SNP calls, with midpoint rooting and collapsed to lineage level. Number of isolates analysed in each lineage (L) is indicated in brackets. **b** Heat map indicating the frequency of each IS6110 insertion site (columns) detected in each lineage (rows), according to inset legend (i.e., a black cell indicates that the IS insertion site was detected in all isolates of the given lineage, while white cell indicates that the insertion site was not detected at all in that lineage). **c** Boxplots show number of IS6110 insertions detected per genome, for each lineage. Black line, median; boxes, interquartile range; whiskers, minimum and maximum values. **d** Histogram of IS6110 insertions in 1000 bp windows along the *M. tuberculosis* chromosome. Dashed line, threshold for defining insertion hotspots

several clusters of insertions detected by ISMapper. Many of these clusters comprised multiple independent insertions into PE/PPE genes (which are surface-associated and interact with the host immune system), as well as the membrane associated proteins *mmpS1* and *mmpL12*. There was substantial clustering of IS6110 insertions interrupting genes encoding the CRISPR machinery, which is involved in immunity to bacteriophage and other foreign DNA. Further, all three phospholipase genes, which are involved in virulence by inducing cell death in macrophages [41] and are encoded by the *plcABC* operon, contained multiple IS6110 insertions detected by ISMapper. This locus is a known hotspot for IS6110 insertions and has been shown to mediate deletions of segments of this region [42]. IS6110 insertions upstream of *phoP*, which have been associated with upregulation and enhanced virulence in *M. tuberculosis* [9], were identified in multiple lineages (1 insertion in 6 lineage 2 genomes; singular insertions in one genome

each in lineage 3 and 5) and may be indicative of positive selection for enhanced *phoP* expression and virulence. These findings from ISMapper analysis are consistent with those reported from PCR-based screens of smaller sets of isolates, but provide a more comprehensive picture of IS dynamics in *M. tuberculosis* that could be extended to much larger genomic data sets and other important pathogens.

Conclusions

ISMapper is a lightweight and reliable tool for the detection of IS insertion sites in bacterial genomes using high-throughput short-read sequencing data, which is now ubiquitous in microbial research and clinical investigations. ISMapper performed well on real and simulated data from 32 different ISs and 13 bacterial species, detecting all but the most complex instances involving multiple neighbouring IS insertions or other repeated sequences. ISMapper was able to detect antimicrobial

resistance-associated ISAbal insertions in *A. baumannii*, with all sites detected by the program being subsequently confirmed by PCR. Compared to other tools such as *breseq* and TIE, ISMapper is ideal for detecting new positions for known ISs in bacterial genomes. In addition, ISMapper was able to rapidly produce a wealth of data on IS6110 insertions in *M. tuberculosis*, allowing quick identification of lineage-specific insertions and specific regions enriched for insertions that may be functionally significant.

Availability and requirements

- **Project name:** ISMapper
- **Project home page:** https://github.com/jhawkey/IS_mapper
- **Programming language:** Python v2.7.5
- **Operating system(s):** platform independent, requires Python 2.7 and dependencies
- **Other requirements:** BioPython v1.63, BWA v0.7.12, SAMtools v1.1, Bedtools v2.20.1, BLAST+ v2.2.28, Samblaster v0.1.21
- **License:** Modified BSD

Competing interests

The authors declare that they have no competing interests.

Authors' contributions

JH developed the code, analysed data and wrote the paper. KEH conceived the study and helped to draft the manuscript. HBJ participated in design and coordination of the study and contributed to data interpretation. RRW and DJE developed code. MH performed PCR and sequence analysis. RMH provided sequence data and isolates for validation and contributed to data interpretation. All authors read and approved the final manuscript.

Acknowledgements

This work was supported by the National Health and Medical Research Council of Australia (Fellowship #1061409 to KEH; Project Grant #1043830 to KEH and RMH) and the Victorian Life Sciences Computation Initiative (VLSI, #VR0082).

Author details

¹Department of Biochemistry and Molecular Biology, Bio21 Molecular Science and Biotechnology Institute, The University of Melbourne, Parkville, VIC 3010, Australia. ²Faculty of Veterinary and Agricultural Science, The University of Melbourne, Parkville, VIC 3010, Australia. ³School of Molecular Bioscience, The University of Sydney, Sydney 2006, Australia.

Received: 9 March 2015 Accepted: 18 August 2015

Published online: 03 September 2015

References

1. Siguier P, Perochon J, Lestrade L, Mahillon J, Chandler M. ISfinder: the reference centre for bacterial insertion sequences. *Nucleic Acids Res.* 2006;34(Database issue):D32–6.
2. Mahillon J, Chandler M. Insertion sequences. *Microbiol Mol Biol Rev.* 1998;62:725–74.
3. Siguier P, Goubeyre E, Chandler M. Bacterial insertion sequences: their genomic impact and diversity. *FEMS Microbiol Rev.* 2014;38:865–91.
4. Oliver A, Vallé M, Chaslus-Dancla E, Cloeckaert A. Overexpression of the multidrug efflux operon *acrEF* by insertional activation with IS1 or IS10 elements in *Salmonella enterica* serovar typhimurium DT204 *acrB* mutants selected with fluoroquinolones. *Antimicrob Agents Chemother.* 2005;49:289–301.
5. Jellen-Ritter AS, Kern WV. Enhanced expression of the multidrug efflux pumps AcrAB and AcrEF associated with insertion element transposition in *Escherichia coli* mutants selected with a fluoroquinolone. *Antimicrob Agents Chemother.* 2001;45:1467–72.
6. Hamidian M, Hall RM. ISAbal targets a specific position upstream of the intrinsic *ampC* gene of *Acinetobacter baumannii* leading to cephalosporin resistance. *J Antimicrob Chemother.* 2013;68:2682–3.
7. Hamidian M, Hancock DP, Hall RM. Horizontal transfer of an ISAbal25-activated *ampC* gene between *Acinetobacter baumannii* strains leading to cephalosporin resistance. *J Antimicrob Chemother.* 2013;68:244–5.
8. Hamidian M, Hall RM. Tn6168, a transposon carrying an ISAbal-activated *ampC* gene and conferring cephalosporin resistance in *Acinetobacter baumannii*. *J Antimicrob Chemother.* 2014;69:77–80.
9. Soto CY, Menéndez MC, Pérez E, Samper S, Gómez AB, García MJ, et al. IS6110 mediates increased transcription of the *phoP* virulence gene in a multidrug-resistant clinical isolate responsible for tuberculosis outbreaks. *J Clin Microbiol.* 2004;42:212–9.
10. Uria MJ, Zhang Q, Li Y, Chan A, Exley RM, Gollan B, et al. A generic mechanism in *Neisseria meningitidis* for enhanced resistance against bactericidal antibodies. *J Exp Med.* 2008;205:1423–34.
11. Van Der Ploeg J, Willemsen M, Van Hall G, Janssen DB. Adaptation of *Xanthobacter autotrophicus* GJ10 to bromoacetate due to activation and mobilization of the haloacetate dehalogenase gene by insertion element IS1247. *J Bacteriol.* 1995;177:1348–56.
12. Aronson BD, Levinthal M, Somerville RL. Activation of a cryptic pathway for threonine metabolism via specific IS3-mediated alteration of promoter structure in *Escherichia coli*. *J Bacteriol.* 1989;171:5503–11.
13. Soria G, Barbé J, Gibert I. Molecular fingerprinting of *Salmonella typhimurium* by IS200-typing as a tool for epidemiological and evolutionary studies. *Microbiologia.* 1994;10:57–68.
14. Das S, Paramasivan CN, Lowrie DB, Prabhakar R, Narayanan PR. IS6110 restriction fragment length polymorphism typing of clinical isolates of *Mycobacterium tuberculosis* from patients with pulmonary tuberculosis in Madras, South India. *Tuber Lung Dis.* 1995;76:550–4.
15. Bik EM, Gouw RD, Mooi FR. DNA fingerprinting of *Vibrio cholerae* strains with a novel insertion sequence element: a tool to identify epidemic strains. *J Clin Microbiol.* 1996;34:1453–61.
16. Adams MD, Chan ER, Molyneux ND, Bonomo RA. Genomewide analysis of divergence of antibiotic resistance determinants in closely related isolates of *Acinetobacter baumannii*. *Antimicrob Agents Chemother.* 2010;54:3569–77.
17. Suzuki M, Matsumoto M, Hata M, Takahashi M, Sakae K. Development of a rapid PCR method using the insertion sequence IS1203 for genotyping Shiga toxin-producing *Escherichia coli* O157. *J Clin Microbiol.* 2004;42:5462–6.
18. Doig KD, Holt KE, Fyfe JA, Lavender CJ, Eddyani M, Portaels F, et al. On the origin of *Mycobacterium ulcerans*, the causative agent of Buruli ulcer. *BMC Genomics.* 2012;13:258.
19. Doughty EL, Sergeant MJ, Adetifa I, Antonio M, Pallen MJ. Culture-independent detection and characterisation of *Mycobacterium tuberculosis* and *M. africanum* in sputum samples using shotgun metagenomics on a benchtop sequencer. *PeerJ.* 2014;2:e585.
20. Rizk G, Gouin A, Chikhi R, Lemaitre C. MindTheGap: integrated detection and assembly of short and long insertions. *Bioinformatics.* 2014;30:3451–7.
21. Chen K, Wallis JW, McLellan MD, Larson DE, Kalicki JM, Pohl CS, et al. BreakDancer: an algorithm for high-resolution mapping of genomic structural variation. *Nat Methods.* 2009;6:677–81.
22. Thung DT, de Ligt J, Vissers LE, Steehouwer M, Kroon M, de Vries P, et al. Mobster: accurate detection of mobile element insertions in next generation sequencing data. *Genome Biol.* 2014;15:488.
23. Robb SMC, Lu L, Valencia E, Burnette JM, Okumoto Y, Wessler SR, et al. The use of RelocaTE and unassembled short reads to produce high-resolution snapshots of transposable element generated diversity in rice. *G3.* 2013;3:949–57.
24. Keane TM, Wong K, Adams DJ. RetroSeq: transposable element discovery from next-generation sequencing data. *Bioinformatics.* 2013;29:389–90.
25. Nakagome M, Solovieva E, Takahashi A, Yasue H, Hirochika H, Miyao A. Transposon Insertion Finder (TIF): a novel program for detection of de novo transpositions of transposable elements. *BMC Bioinformatics.* 2014;15:71.
26. Barrick JE, Colburn G, Deatherage DE, Traverse CC, Strand MD, Borges JJ, et al. Identifying structural variation in haploid microbial genomes from short-read resequencing data using breseq. *BMC Genomics.* 2014;15:1039.

27. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*. 2009;25:1754–60.
28. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The sequence alignment/Map format and SAMtools. *Bioinformatics*. 2009;25:2078–9.
29. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics*. 2010;26:841–2.
30. Faust GG, Hall IM. SAMBLASTER: fast duplicate marking and structural variant read extraction. *Bioinformatics*. 2014;30:2503–5.
31. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+: architecture and applications. *BMC Bioinformatics*. 2009;10:421.
32. Holt KE, Parkhill J, Mazzoni CJ, Roumagnac P, Weill F-X, Goodhead I, et al. High-throughput sequencing provides insights into genome variation and evolution in *Salmonella* Typhi. *Nat Genet*. 2008;40:987–93.
33. Inouye M, Dashnow H, Raven L-A, Schultz MB, Pope BJ, Tomita T, et al. SRST2: rapid genomic surveillance for public health and hospital microbiology labs. *Genome Med*. 2014;6:90.
34. Gupta SK, Padmanabhan BR, Diene SM, Lopez-Rojas R, Kempf M, Landraud L, Rolain J-M: ARG-ANNOT, a new bioinformatic tool to discover antibiotic resistance genes in bacterial genomes. *Antimicrob Agents Chemother* 2014;58:212–20.
35. McEvoy CRE, Falmer AA, van Pittius NCG, Victor TC, van Helden PD, Warren RM. The role of IS6110 in the evolution of *Mycobacterium tuberculosis*. *Tuberculosis*. 2007;87:393–404.
36. van Embden JD, Cave MD, Crawford JT, Dale JW, Eisenach KD, Gicquel B, et al. Strain identification of *Mycobacterium tuberculosis* by DNA fingerprinting: recommendations for a standardized methodology. *J Clin Microbiol*. 1993;31:406–9.
37. Beggs ML, Eisenach KD, Cave MD. Mapping of IS6110 insertion sites in two epidemic strains of *Mycobacterium tuberculosis*. *J Clin Microbiol*. 2000;38:2923–8.
38. Alonso H, Aguilo JI, Samper S, Caminero JA, Campos-Herrero MI, Gicquel B, et al. Deciphering the role of IS6110 in a highly transmissible *Mycobacterium tuberculosis* Beijing strain, GC1237. *Tuberculosis*. 2011;91:117–26.
39. Comas I, Coscolla M, Luo T, Borrell S, Holt KE, Kato-Maeda M, et al. Out-of-Africa migration and Neolithic coexpansion of *Mycobacterium tuberculosis* with modern humans. *Nat Genet*. 2013;45:1176–82.
40. Holt KE, Thieu Nga TV, Thanh DP, Vinh H, Kim DW, Vu Tra MP, et al. Tracking the establishment of local endemic populations of an emergent enteric pathogen. *Proc Natl Acad Sci U S A*. 2013;110:17522–7.
41. Assis PA, Espíndola MS, Paula-Silva FW, Rios WM, Pereira PA, Leão SC, et al. *Mycobacterium tuberculosis* expressing phospholipase C subverts PGE₂ synthesis and induces necrosis and alveolar macrophages. *BMC Microbiol*. 2014;14:128.
42. Vera-Cabrera L, Hernández-Vera MA, Welsh O, Johnson WM, Castro-Garza J. Phospholipase region of *Mycobacterium tuberculosis* is a preferential locus for IS6110 transposition. *J Clin Microbiol*. 2001;39:3499–504.

Submit your next manuscript to BioMed Central and take full advantage of:

- Convenient online submission
- Thorough peer review
- No space constraints or color figure charges
- Immediate publication on acceptance
- Inclusion in PubMed, CAS, Scopus and Google Scholar
- Research which is freely available for redistribution

Submit your manuscript at
www.biomedcentral.com/submit

